

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2020-27168

(P2020-27168A)

(43) 公開日 令和2年2月20日(2020.2.20)

(51) Int.Cl.  
G10L 25/30 (2013.01)

F I  
G I O L 25/30

テーマコード (参考)

審査請求 未請求 請求項の数 11 O L (全 19 頁)

<p>(21) 出願番号 特願2018-151611 (P2018-151611)</p> <p>(22) 出願日 平成30年8月10日 (2018. 8. 10)</p>	<p>(71) 出願人 504202472 大学共同利用機関法人情報・システム研究機構 東京都立川市緑町10番3号</p> <p>(74) 代理人 100107766 弁理士 伊東 忠重</p> <p>(74) 代理人 100070150 弁理士 伊東 忠彦</p> <p>(72) 発明者 ヒュウ ティ ルオン 東京都千代田区一ツ橋二丁目1番2号 大学共同利用機関法人情報・システム研究機構 国立情報学研究所内</p> <p>(72) 発明者 山岸 順一 東京都千代田区一ツ橋二丁目1番2号 大学共同利用機関法人情報・システム研究機構 国立情報学研究所内</p>
--	--

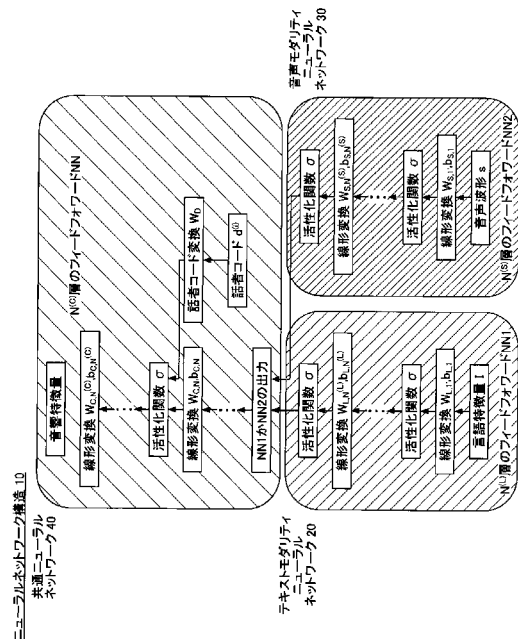
(54) 【発明の名称】 学習装置、学習方法、音声合成装置、音声合成方法及びプログラム

(57) 【要約】 (修正有)

【課題】教師有り適応と教師なし適応との何れのケースにも対応可能な、ニューラルネットワーク構造を持った未知話者のための音声合成技術を提供する。

【解決手段】テキストデータを第1のベクトルに変換するテキストモダリティニューラルネットワーク(テキストモダリティNN)と、音声波形データを第2のベクトルに変換する音声モダリティNNと、第1又は第2のベクトルから話者空間上の話者コードベクトルに対応する音響特徴量を生成する共通NNを備え、テキストデータと音響特徴量で構成される第1の訓練データによりテキストモダリティNN及び共通NNを学習し、音声波形データと音響特徴量で構成される第2の訓練データにより音声モダリティNN及び共通NNを学習し、所与の話者の第3の訓練データに応じて、テキストモダリティNN及び共通NNと、音声モダリティNN及び共通NNとを選択的に利用して、話者に対する話者コードベクトルを推定する。

【選択図】 図1



## 【特許請求の範囲】

## 【請求項 1】

メモリと、  
プロセッサと、  
を有する学習装置であって、  
前記メモリは、  
テキストデータを第 1 のベクトルに変換するテキストモダリティニューラルネットワークと、  
音声波形データを第 2 のベクトルに変換する音声モダリティニューラルネットワークと、

10

前記テキストモダリティニューラルネットワーク及び前記音声モダリティニューラルネットワークに接続され、前記第 1 のベクトル又は前記第 2 のベクトルから話者空間上の話者コードベクトルに対応する音響特徴量を生成する共通ニューラルネットワークとを格納し、

前記プロセッサは、

テキストデータと音響特徴量とから構成される第 1 の訓練データによって前記テキストモダリティニューラルネットワーク及び前記共通ニューラルネットワークを学習し、

音声波形データと音響特徴量とから構成される第 2 の訓練データによって前記音声モダリティニューラルネットワーク及び前記共通ニューラルネットワークを学習し、

所与の話者の第 3 の訓練データに応じて、前記テキストモダリティニューラルネットワーク及び前記共通ニューラルネットワークと、前記音声モダリティニューラルネットワーク及び前記共通ニューラルネットワークとを選択的に利用して、前記所与の話者に対する前記話者コードベクトルを推定する学習装置。

20

## 【請求項 2】

前記プロセッサは、

前記第 1 の訓練データのテキストデータを前記テキストモダリティニューラルネットワークに入力し、前記テキストモダリティニューラルネットワークから取得した第 1 のベクトルを前記共通ニューラルネットワークに入力し、前記共通ニューラルネットワークから取得した音響特徴量と前記第 1 の訓練データの音響特徴量との間の第 1 の誤差を算出し、

前記第 2 の訓練データの音声波形データを前記音声モダリティニューラルネットワークに入力し、前記音声モダリティニューラルネットワークから取得した第 2 のベクトルを前記共通ニューラルネットワークに入力し、前記共通ニューラルネットワークから取得した音響特徴量と前記第 2 の訓練データの音響特徴量との間の第 2 の誤差を算出し、

30

前記第 1 の誤差と前記第 2 の誤差との加重和に基づき、前記テキストモダリティニューラルネットワーク、前記音声モダリティニューラルネットワーク及び前記共通ニューラルネットワークを学習する、請求項 1 記載の学習装置。

## 【請求項 3】

前記プロセッサは、

前記第 1 の訓練データのテキストデータを前記テキストモダリティニューラルネットワークに入力し、前記テキストモダリティニューラルネットワークから取得した第 1 のベクトルを前記共通ニューラルネットワークに入力し、前記共通ニューラルネットワークから取得した音響特徴量と前記第 1 の訓練データの音響特徴量との間の第 1 の誤差を算出し、

40

前記第 2 の訓練データの音声波形データを前記音声モダリティニューラルネットワークに入力し、前記音声モダリティニューラルネットワークから取得した第 2 のベクトルを前記共通ニューラルネットワークに入力し、前記共通ニューラルネットワークの一部のレイヤから構成されるサブニューラルネットワークから第 3 のベクトルを取得し、前記共通ニューラルネットワークに入力された第 1 のベクトルに対して前記サブニューラルネットワークから第 4 のベクトルを取得し、前記第 3 のベクトルと前記第 4 のベクトルとの間の距離に基づき第 3 の誤差を算出し、

前記第 1 の誤差と前記第 3 の誤差との加重和に基づき、前記テキストモダリティニュー

50

ラルネットワーク、前記音声モダリティニューラルネットワーク及び前記共通ニューラルネットワークを学習する、請求項 1 又は 2 記載の学習装置。

【請求項 4】

前記プロセッサは、

前記第 3 の訓練データがテキストデータと音響特徴量とから構成される場合、前記テキストデータを前記テキストモダリティニューラルネットワークに入力し、前記テキストモダリティニューラルネットワークから取得した第 1 のベクトルを前記共通ニューラルネットワークに入力し、前記共通ニューラルネットワークから取得した音響特徴量と前記第 3 の訓練データの音響特徴量との間の第 4 の誤差に基づき前記所与の話者の話者コードベクトルを決定する、請求項 1 乃至 3 何れか一項記載の学習装置。

10

【請求項 5】

前記プロセッサは、

前記第 3 の訓練データが音声波形データと音響特徴量とから構成される場合、前記音声波形データを前記音声モダリティニューラルネットワークに入力し、前記音声モダリティニューラルネットワークから取得した第 2 のベクトルを前記共通ニューラルネットワークに入力し、前記共通ニューラルネットワークから取得した音響特徴量と前記第 3 の訓練データの音響特徴量との間の第 5 の誤差に基づき前記所与の話者の話者コードベクトルを決定する、請求項 1 乃至 4 何れか一項記載の学習装置。

【請求項 6】

メモリと、

プロセッサと、

を有する音声合成装置であって、

前記メモリは、

学習済みのテキストモダリティニューラルネットワークと、

所与の話者に対して学習済みの共通ニューラルネットワークと、

を格納し、

前記プロセッサは、テキストデータを取得すると、前記格納されているテキストモダリティニューラルネットワーク及び共通ニューラルネットワークによって、前記テキストデータから前記所与の話者に対応する音響特徴量を生成する音声合成装置。

20

【請求項 7】

テキストデータを取得し、前記所与の話者に対応して前記テキストデータから生成された音響特徴量を再生する入出力インタフェースを更に有する、請求項 6 記載の音声合成装置。

30

【請求項 8】

メモリとプロセッサとを有するコンピュータによって実現される学習方法であって、

前記メモリは、

テキストデータを第 1 のベクトルに変換するテキストモダリティニューラルネットワークと、

音声波形データを第 2 のベクトルに変換する音声モダリティニューラルネットワークと

、

前記テキストモダリティニューラルネットワーク及び前記音声モダリティニューラルネットワークに接続され、前記第 1 のベクトル又は前記第 2 のベクトルから話者空間上の話者コードベクトルに対応する音響特徴量を生成する共通ニューラルネットワークとを格納し、

40

前記プロセッサが、テキストデータと音響特徴量とから構成される第 1 の訓練データによって前記テキストモダリティニューラルネットワーク及び前記共通ニューラルネットワークを学習するステップと、

前記プロセッサが、音声波形データと音響特徴量とから構成される第 2 の訓練データによって前記音声モダリティニューラルネットワーク及び前記共通ニューラルネットワークを学習するステップと、

50

前記プロセッサが、所与の話者の第3の訓練データに応じて、前記テキストモダリティニューラルネットワーク及び前記共通ニューラルネットワークと、前記音声モダリティニューラルネットワーク及び前記共通ニューラルネットワークとを選択的に利用して、前記所与の話者に対する前記話者コードベクトルを推定するステップと、  
を有する学習方法。

【請求項9】

メモリとプロセッサとを有するコンピュータによって実現される音声合成方法であって

前記メモリは、  
学習済みのテキストモダリティニューラルネットワークと、  
所与の話者に対して学習済みの共通ニューラルネットワークと、  
を格納し、

前記プロセッサが、テキストデータを取得すると、前記格納されているテキストモダリティニューラルネットワーク及び共通ニューラルネットワークによって、前記テキストデータから前記所与の話者に対応する音響特徴量を生成するステップを有する音声合成方法。

【請求項10】

テキストデータを第1のベクトルに変換するテキストモダリティニューラルネットワークと、音声波形データを第2のベクトルに変換する音声モダリティニューラルネットワークと、前記テキストモダリティニューラルネットワーク及び前記音声モダリティニューラルネットワークに接続され、前記第1のベクトル又は前記第2のベクトルから話者空間上の話者コードベクトルに対応する音響特徴量を生成する共通ニューラルネットワークとを格納したメモリに接続されるプロセッサに、

テキストデータと音響特徴量とから構成される第1の訓練データによって前記テキストモダリティニューラルネットワーク及び前記共通ニューラルネットワークを学習させ、

音声波形データと音響特徴量とから構成される第2の訓練データによって前記音声モダリティニューラルネットワーク及び前記共通ニューラルネットワークを学習させ、

所与の話者の第3の訓練データに応じて、前記テキストモダリティニューラルネットワーク及び前記共通ニューラルネットワークと、前記音声モダリティニューラルネットワーク及び前記共通ニューラルネットワークとを選択的に利用して、前記所与の話者に対する前記話者コードベクトルを推定させるプログラム。

【請求項11】

学習済みのテキストモダリティニューラルネットワークと、所与の話者に対して学習済みの共通ニューラルネットワークとを格納したメモリに接続されるプロセッサに、

テキストデータを取得すると、前記格納されているテキストモダリティニューラルネットワーク及び共通ニューラルネットワークによって、前記テキストデータから前記所与の話者に対応する音響特徴量を生成させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般に音声合成技術に関し、より詳細には、ニューラルネットワークを利用した未知話者に対する話者適応技術に関する。

【背景技術】

【0002】

近年のディープラーニングの進展によって、ニューラルネットワークを利用した音声合成システムの研究開発が進められている。

【0003】

音声合成システムの一例として、特定話者のための音声合成システムがある。特定話者のための音声合成システムによると、特定話者の音声データとテキストデータとのペアを訓練データとして利用することによって、テキストデータを当該話者に対応する音声デー

10

20

30

40

50

タに変換するニューラルネットワークが学習され、学習済みのニューラルネットワークを利用して、入力されたテキストデータが当該特定話者の音声によって再生される。

【0004】

他の例として、複数話者のための音声合成システムがある。複数話者のための音声合成システムによると、複数話者の音声データとテキストデータとのペアを訓練データとして利用することによって、テキストデータを複数話者の何れか指定された話者に対応する音声データに変換するニューラルネットワークが学習され、学習済みのニューラルネットワークを利用して、入力されたテキストデータが当該指定された話者の音声によって再生される。

【0005】

更なる他の例として、未知話者のための音声合成システムがある。典型的には、上述した複数話者のための音声合成システムに基づき、未知話者の音声データ及び/又はテキストデータを訓練データとして利用することによって、テキストデータを当該未知話者に対応する音声データに変換するニューラルネットワークが学習される。学習済みのニューラルネットワークを利用して、入力されたテキストデータが当該未知話者の音声によって再生される。

10

【0006】

未知話者のための音声合成システムとして、未知話者の音声データとテキストデータとのペアを訓練データとして利用するもの（教師有り適応と呼ばれる）と、未知話者の音声データのみを訓練データとして利用するもの（教師なし適応と呼ばれる）とがある。

20

【先行技術文献】

【非特許文献】

【0007】

【非特許文献1】"Neural Voice Cloning with a Few Samples", Sercan O. Arik, et. al., arXiv: 1802.06006, Mar. 20, 2018.

【非特許文献2】"Fitting New Speakers Based on a Short Untranscribed Sample", Eliya Nachmani, et. al., arXiv: 1802.06984, Feb. 20, 2018.

【発明の概要】

【発明が解決しようとする課題】

【0008】

従来技術によると、教師有り適応に基づく未知話者のための音声合成システムと、教師なし適応に基づく未知話者のための音声合成システムとは、それぞれ独立に設計されており、教師有り適応と教師なし適応との双方に対応可能な音声合成システムは現状存在しない。従って、教師有り適応と教師なし適応との何れのケースにも対応可能なニューラルネットワーク構造を備えた未知話者のための音声合成システムが望まれる。

30

【0009】

上述した問題点を鑑み、本発明の課題は、教師有り適応と教師なし適応との何れのケースにも対応可能なニューラルネットワーク構造を利用した未知話者のための音声合成技術を提供することである。

【課題を解決するための手段】

40

【0010】

上記課題を解決するため、本発明の一態様は、メモリと、プロセッサとを有する学習装置であって、前記メモリは、テキストデータを第1のベクトルに変換するテキストモダリティニューラルネットワークと、音声波形データを第2のベクトルに変換する音声モダリティニューラルネットワークと、前記テキストモダリティニューラルネットワーク及び前記音声モダリティニューラルネットワークに接続され、前記第1のベクトル又は前記第2のベクトルから話者空間上の話者コードベクトルに対応する音響特徴量を生成する共通ニューラルネットワークとを格納し、前記プロセッサは、テキストデータと音響特徴量とから構成される第1の訓練データによって前記テキストモダリティニューラルネットワーク及び前記共通ニューラルネットワークを学習し、音声波形データと音響特徴量とから構成

50

される第2の訓練データによって前記音声モダリティニューラルネットワーク及び前記共通ニューラルネットワークを学習し、所与の話者の第3の訓練データに応じて、前記テキストモダリティニューラルネットワーク及び前記共通ニューラルネットワークと、前記音声モダリティニューラルネットワーク及び前記共通ニューラルネットワークとを選択的に利用して、前記所与の話者に対する前記話者コードベクトルを推定する学習装置に関する。

【発明の効果】

【0011】

本発明によると、教師有り適応と教師なし適応との何れのケースにも対応可能なニューラルネットワーク構造を利用した未知話者のための音声合成技術を提供することができる。

10

【図面の簡単な説明】

【0012】

【図1】本発明の一実施例によるニューラルネットワーク構造の概略図である。

【図2】本発明の一実施例による学習装置及び音声合成装置のハードウェア構成を示すブロック図である。

【図3】本発明の一実施例による学習処理を示す概略図である。

【図4】本発明の一実施例による学習処理を示すフローチャートである。

【図5】本発明の他の実施例による学習処理を示す概略図である。

【図6】本発明の他の実施例による学習処理を示すフローチャートである。

20

【図7】本発明の一実施例による未知話者適応処理を示す概略図である。

【図8】本発明の一実施例による未知話者適応処理を示すフローチャートである。

【図9】本発明の一実施例による音声合成処理を示す概略図である。

【図10】本発明の一実施例による音声合成処理を示すフローチャートである。

【図11】本発明の各種実施例による学習処理の実験結果を示す図である。

【発明を実施するための形態】

【0013】

以下の実施例では、教師有り適応と教師なし適応との何れのケースにも対応可能なニューラルネットワークを学習する学習装置100と、当該ニューラルネットワークを利用した未知話者のための音声合成装置200とが開示される。

30

[概略]

後述される実施例を概略すると、学習装置100は、テキストデータをベクトルに変換するテキストモダリティニューラルネットワーク20、音声波形データをベクトルに変換する音声モダリティニューラルネットワーク30、及びテキストモダリティニューラルネットワーク20及び音声モダリティニューラルネットワーク30から出力されたベクトルから、話者空間上の所与の未知話者を示す話者コードベクトル（潜在変数）に対応する音響特徴量を生成する共通ニューラルネットワーク40を学習する。

【0014】

まず、テキストモダリティニューラルネットワーク20、音声モダリティニューラルネットワーク30及び共通ニューラルネットワーク40から構成されるニューラルネットワーク構造10に対する学習処理において、学習装置100は、テキストデータと音響特徴量とのペアから構成される訓練データに対して、テキストデータをテキストモダリティニューラルネットワーク20に入力し、テキストモダリティニューラルネットワーク20から出力されたベクトルを共通ニューラルネットワーク40に入力する。一方、学習装置100は、音声波形データと音響特徴量とのペアから構成される訓練データに対して、音声波形データを音声モダリティニューラルネットワーク30に入力し、音声モダリティニューラルネットワーク30から取得されたベクトルを共通ニューラルネットワーク40に入力する。そして、以下の実施例において詳細に説明されるように、学習装置100は、共通ニューラルネットワーク40から出力された音響特徴量と訓練データの音響特徴量とに基づき、テキストモダリティニューラルネットワーク20、音声モダリティニューラルネ

40

50

ットワーク 30 及び共通ニューラルネットワーク 40 を学習する。

【0015】

次に、未知話者適応処理において、学習装置 100 は、上述した学習済みのテキストモダリティニューラルネットワーク 20、音声モダリティニューラルネットワーク 30 及び共通ニューラルネットワーク 40 を利用して、話者空間上の未知話者の位置を示す話者コードベクトルを推定する。すなわち、所与の話者の訓練データが与えられると、学習装置 100 は、当該訓練データがテキスト付きの音声データであるか、あるいは、音声データのみであるかに応じて、テキストモダリティニューラルネットワーク 20 及び共通ニューラルネットワーク 40 と、音声モダリティニューラルネットワーク 30 及び共通ニューラルネットワーク 40 とを選択的に利用して、共通ニューラルネットワーク 40 の話者空間上の当該話者を示す話者コードベクトル（潜在変数）を推定し、推定した潜在変数が埋め込まれた話者毎の共通ニューラルネットワーク 40 を生成する。

10

【0016】

音声合成装置 200 は、このようにして学習装置 100 によって未知話者毎に学習されたニューラルネットワーク構造 10 における学習済みのテキストモダリティニューラルネットワーク 20 及び共通ニューラルネットワーク 40 を利用して、所与のテキストデータから当該未知話者に対応する音声データを生成する。

[ニューラルネットワーク構造]

まず、図 1 を参照して、本発明の一実施例によるニューラルネットワーク構造 10 を説明する。図 1 は、本発明の一実施例によるニューラルネットワーク構造 10 の概略図である。

20

【0017】

図 1 に示されるように、ニューラルネットワーク構造 10 は、テキストモダリティニューラルネットワーク 20、音声モダリティニューラルネットワーク 30 及び共通ニューラルネットワーク 40 を有し、テキストモダリティニューラルネットワーク 20 及び音声モダリティニューラルネットワーク 30 はそれぞれ、共通ニューラルネットワーク 40 に接続される。

【0018】

テキストモダリティニューラルネットワーク 20 は、入力されたテキストデータ（例えば、言語特徴量）を共通ニューラルネットワーク 40 への入力用のベクトルに変換する何れかのレイヤ構成を有するニューラルネットワークである。図示された実施例では、テキストモダリティニューラルネットワーク 20 は、 $N^{(L)}$  層のフィードフォワードニューラルネットワークであり、テキストデータのベクトル  $l$  を入力層において取得し、取得したベクトル  $l$  を隠れ層にわたす。 $N^{(L)}$  個の隠れ層はそれぞれ、前段のレイヤからわたされたベクトルを行列  $W_L$  及びバイアスベクトル  $b_L$  によって線形変換し、変換されたベクトルを活性化関数（例えば、シグモイド関数）に入力し、活性化関数から出力されたベクトルを後段のレイヤにわたす。出力層は、前段の隠れ層からわたされたベクトルを共通ニューラルネットワーク 40 の入力層にわたす。

30

【0019】

形式的には、テキストデータのベクトル  $l$  が与えられると、第 1 の隠れ層は、

40

$$h_1 = (W_{L,1} l + b_{L,1})$$

によってベクトル  $h_1$  を出力する。以下同様にして、各隠れ層は同様の変換処理を実行し、第  $N^{(L)}$  の隠れ層は、前段の隠れ層からベクトル  $h_{N^{(L)}-1}$  が与えられると、

$$h_{N^{(L)}} = (W_{L,N^{(L)}} h_{N^{(L)}-1} + b_{L,N^{(L)}})$$

によってベクトル  $h_{N^{(L)}}$  を出力し、出力層にわたす。当該ベクトル及び行列は、後述される学習処理において学習される。

【0020】

音声モダリティニューラルネットワーク 30 は、入力された音声データ（例えば、音声波形）を共通ニューラルネットワーク 40 への入力用のベクトルに変換する何れかのレイヤ構成を有するニューラルネットワークである。図示された実施例では、音声モダリティ

50

ニューラルネットワーク 30 は、 $N^{(s)}$  層のフィードフォワードニューラルネットワークであり、音声データのベクトル  $s$  を入力層において取得し、取得したベクトル  $s$  を隠れ層にわたす。 $N^{(s)}$  個の隠れ層はそれぞれ、前段のレイヤからわたされたベクトルを行列  $W_s$  及びバイアスベクトル  $b_s$  によって線形変換し、変換されたベクトルを活性化関数（例えば、シグモイド関数）に入力し、活性化関数 から出力されたベクトルを後段のレイヤにわたす。出力層は、前段の隠れ層からわたされたベクトルを共通ニューラルネットワーク 40 の入力層にわたす。なお、各隠れ層における具体的な処理は、上述したテキストモダリティニューラルネットワーク 20 のものと同様であり、重複する説明は省く。

【0021】

共通ニューラルネットワーク 40 は、テキストモダリティニューラルネットワーク 20 及び音声モダリティニューラルネットワーク 30 からわたされたベクトルを音響特徴量に変換する何れかのレイヤ構成を有するニューラルネットワークである。図示された実施例では、共通ニューラルネットワーク 40 は、 $N^{(c)}$  層のフィードフォワードニューラルネットワークであり、テキストモダリティニューラルネットワーク 20 及び音声モダリティニューラルネットワーク 30 から入力されたベクトルを入力層において取得し、取得したベクトルを隠れ層にわたす。 $N^{(c)}$  個の隠れ層はそれぞれ、前段のレイヤからわたされたベクトルを行列  $W_c$  及びバイアスベクトル  $b_c$  によって線形変換し、変換されたベクトルを活性化関数（例えば、シグモイド関数）に入力し、活性化関数 から出力されたベクトルを後段のレイヤにわたす。出力層は、前段の隠れ層からわたされた音響特徴量を示すベクトルを出力する。

【0022】

また、共通ニューラルネットワーク 40 は更に、後述される未知話者適応処理によって推定された所与の話者を示す話者コードベクトル（潜在変数）を含む。換言すると、共通ニューラルネットワーク 40 は、未知話者適応処理において学習装置 100 によって話者毎に学習される。所与の話者を示す話者空間上の推定された話者コードベクトルが与えられた隠れ層は、前段のレイヤからわたされたベクトルと話者コードベクトルとに対して線形変換を実行し、変換されたベクトルを活性化関数（例えば、シグモイド関数）に入力し、活性化関数 から出力されたベクトルを後段のレイヤにわたす。

【0023】

形式的には、話者コードベクトルが与えられる隠れ層は、前段のレイヤからベクトル  $h_{n-1}$  と話者  $i$  の話者コードベクトル  $d^{(i)}$  とが与えられると、  

$$h_n = (W_{c,n} h_{n-1} + b_{c,n} + W_D d^{(i)})$$
 によってベクトル  $h_n$  を取得する。ここで、 $W_D$  は話者コード用の重み行列である。なお、話者コードベクトルが入力されない各隠れ層における具体的な処理は、上述したテキストモダリティニューラルネットワーク 20 のものと同様であり、重複する説明は省く。

【0024】

なお、図示された実施例では、話者コードベクトルは 1 つの隠れ層にわたされているが、これに限定されるものでなく、共通ニューラルネットワーク 40 のレイヤ構成に応じて複数の隠れ層にわたされてもよい。

[ハードウェア構成]

ここで、学習装置 100 及び音声合成装置 200 は、例えば、図 2 に示されるように、CPU (Central Processing unit)、GPU (Graphics Processing Unit) などのプロセッサ 101、RAM (Random Access Memory)、フラッシュメモリなどのメモリ 102、ハードディスク 103 及び入出力(I/O)インタフェース 104 によるハードウェア構成を有してもよい。

【0025】

プロセッサ 101 は、学習装置 100 及び音声合成装置 200 の各種処理を実行する。

【0026】

メモリ 102 は、学習装置 100 及び音声合成装置 200 における各種データ及びプログラムを格納し、特に作業用データ、実行中のプログラムなどのためのワーキングメモリ

10

20

30

40

50



として機能する。具体的には、メモリ 102 は、ハードディスク 103 からロードされたニューラルネットワーク構造 10 を実現するプログラム、各種処理を実行及び制御するためのプログラムなどを格納し、プロセッサ 101 によるプログラムの実行中にワーキングメモリとして機能する。

【0027】

ハードディスク 103 は、学習装置 100 及び音声合成装置 200 における各種データ及びプログラムを格納する。

【0028】

I/Oインタフェース 104 は、ユーザからの命令、入力データなどを受け付け、出力結果を表示、再生などすると共に、外部装置との間でデータを入出力するためのインタフェースである。例えば、I/Oインタフェース 104 は、USB (Universal Serial Bus)、通信回線、キーボード、マウス、ディスプレイ、マイクロフォン、スピーカなどの各種データを入出力するためのデバイスである。

【0029】

しかしながら、本発明による学習装置 100 及び音声合成装置 200 は、上述したハードウェア構成に限定されず、他の何れか適切なハードウェア構成を有してもよい。例えば、上述した学習装置 100 及び音声合成装置 200 による各種処理の 1 つ以上は、これを実現するよう配線化された処理回路又は電子回路により実現されてもよい。

[ニューラルネットワーク構造の第 1 の学習処理]

次に、図 3 及び 4 を参照して、本発明の一実施例によるニューラルネットワーク構造 10 に対する学習処理を説明する。上述したニューラルネットワーク構造 10 の内部構成から理解されるように、学習装置 100 は、共通ニューラルネットワーク 40 がテキストデータと音声データとの異なるモダリティからの入力を適切に受け付けるようにニューラルネットワーク構造 10 を学習する必要がある。

【0030】

図 3 は、本発明の一実施例による学習処理を示す概略図である。本実施例では、図 3 に示されるように、学習装置 100 は、共通ニューラルネットワーク 40 をテキストモダリティニューラルネットワーク 20 と音声モダリティニューラルネットワーク 30 とに共有させ、2 つの共通ニューラルネットワーク 40 を同時に、すなわち、2 つの共通ニューラルネットワーク 40 におけるパラメータ (例えば、隠れ層の重み行列) が同一となるよう同期的に学習する。

【0031】

具体的には、学習装置 100 は、テキストデータと音響特徴量とのペアから構成される訓練データに対して、当該テキストデータをテキストモダリティニューラルネットワーク 20 に入力し、テキストモダリティニューラルネットワーク 20 から出力されたベクトルを取得する。そして、学習装置 100 は、取得したベクトルを共通ニューラルネットワーク 40 に入力し、共通ニューラルネットワーク 40 から出力された音響特徴量を取得し、取得した音響特徴量と訓練データの音響特徴量との間の誤差 ( $loss_{main}$ ) を算出する。

【0032】

一方、学習装置 100 は、音声波形データと音響特徴量とのペアから構成される訓練データに対して、当該音声波形データを音声モダリティニューラルネットワーク 30 に入力し、音声モダリティニューラルネットワーク 30 から出力されたベクトルを取得する。そして、学習装置 100 は、取得したベクトルを共通ニューラルネットワーク 40 に入力し、共通ニューラルネットワーク 40 から出力された音響特徴量を取得し、取得した音響特徴量と訓練データの音響特徴量との間の誤差 ( $loss_{sub}$ ) を算出する。

【0033】

その後、学習装置 100 は、算出した 2 つの誤差 ( $loss_{main}$ ,  $loss_{sub}$ ) の加重和に基づき、テキストモダリティニューラルネットワーク 20、音声モダリティニューラルネットワーク 30 及び共通ニューラルネットワーク 40 を学習する。例えば、

10

20

30

40

50

学習装置 100 は、

$$loss = loss_{main} + loss_{sub}$$

に従って (  $loss$  は、スカラー値である )、テキストモダリティニューラルネットワーク 20 及び共通ニューラルネットワーク 40 による誤差  $loss_{main}$  と、音声モダリティニューラルネットワーク 30 及び共通ニューラルネットワーク 40 による誤差  $loss_{sub}$  との 2 つの誤差の加重和 (  $loss$  ) を算出してもよい。

【0034】

学習装置 100 は、算出した誤差の加重和 (  $loss$  ) が減少するように、例えば、バックプロパゲーションに従って、共有される 2 つの共通ニューラルネットワーク 40 のパラメータが同一となるように、テキストモダリティニューラルネットワーク 20、音声モダリティニューラルネットワーク 30 及び共通ニューラルネットワーク 40 のパラメータ (例えば、隠れ層の重み行列) を更新する。

10

【0035】

図 4 は、本発明の一実施例による学習処理を示すフローチャートである。当該学習処理は、学習装置 100、具体的には、学習装置 100 のプロセッサ 101 によって実行される。

【0036】

図 4 に示されるように、ステップ S 101 において、学習装置 100 は、訓練データを取得する。例えば、訓練データが複数の話者によるテキスト付きの音声データである場合、学習装置 100 は、前処理として、当該音声データに対応する音声波形データ及び音響特徴量に変換し、訓練データからテキストデータと音響特徴量とのペアと音声波形データと音響特徴量とのペアとを生成してもよい。

20

【0037】

ステップ S 102 において、学習装置 100 は、処理対象の訓練データがテキストデータと音響特徴量とのペアである場合、ステップ S 103 に進み、処理対象の訓練データが音声波形データと音響特徴量とのペアである場合、ステップ S 106 に進む。

【0038】

ステップ S 103 において、学習装置 100 は、訓練データのテキストデータをテキストモダリティニューラルネットワーク 20 に入力し、テキストモダリティニューラルネットワーク 20 から出力されたベクトルを取得する。

30

【0039】

ステップ S 104 において、学習装置 100 は、取得したベクトルを共通ニューラルネットワーク 40 に入力し、共通ニューラルネットワーク 40 から出力された音響特徴量を取得する。

【0040】

ステップ S 105 において、学習装置 100 は、共通ニューラルネットワーク 40 から取得した音響特徴量と訓練データの音響特徴量との誤差 (  $loss_{main}$  ) を算出する。

【0041】

一方、ステップ S 106 において、学習装置 100 は、訓練データの音声波形データを音声モダリティニューラルネットワーク 30 に入力し、音声モダリティニューラルネットワーク 30 から出力されたベクトルを取得する。

40

【0042】

ステップ S 107 において、学習装置 100 は、取得したベクトルを共通ニューラルネットワーク 40 に入力し、共通ニューラルネットワーク 40 から出力された音響特徴量を取得する。

【0043】

ステップ S 108 において、学習装置 100 は、共通ニューラルネットワーク 40 から取得した音響特徴量と訓練データの音響特徴量との誤差 (  $loss_{sub}$  ) を算出する。

【0044】

50

ステップS109において、学習装置100は、ステップS105及びS108において取得した2つの誤差の加重和(LOSS)を計算し、計算した加重和(LOSS)が減少するように、例えば、バックプロパゲーションに従ってテキストモダリティニューラルネットワーク20、音声モダリティニューラルネットワーク30及び共通ニューラルネットワーク40のパラメータ(例えば、隠れ層の重み行列)を更新し、具体的には、共有される2つの共通ニューラルネットワーク40のパラメータが同一のものに更新されるように、2つの共通ニューラルネットワーク40を同期的に学習する。

#### 【0045】

学習装置100は、所定の終了条件を充足するまで、各訓練データに対して上述したステップS101～S109を繰り返す。当該所定の終了条件は、例えば、所定の回数の繰り返しを終了したこと、誤差(LOSS)が所定の閾値以下になったこと、誤差(LOSS)が収束したことなどであってもよい。

#### [ニューラルネットワーク構造の第2の学習処理]

次に、図5及び6を参照して、本発明の他の実施例によるニューラルネットワーク構造10に対する学習処理を説明する。上述したニューラルネットワーク構造10から理解されるように、学習装置100は、共通ニューラルネットワーク40がテキストデータと音声データとの異なるモダリティからの入力を適切に受け付けるようにニューラルネットワーク構造10、特に、共通ニューラルネットワーク40の入力層に近い下層レイヤを学習することが求められる。

#### 【0046】

図5は、本発明の他の実施例による学習処理を示す概略図である。本実施例では、図5に示されるように、学習装置100は、テキストモダリティニューラルネットワーク20及び音声モダリティニューラルネットワーク30から入力された各ベクトルに対して、共通ニューラルネットワーク40における一部の隠れ層(例えば、入力層から所定番目の隠れ層)から出力される各ベクトルの間の距離を損失又はペナルティ(LOSS<sub>sub</sub>)として利用し、上述したテキストモダリティニューラルネットワーク20及び共通ニューラルネットワーク40における誤差(LOSS<sub>main</sub>)と、一部の隠れ層から出力されるベクトル間の距離(LOSS<sub>sub</sub>)との加重和に基づき、テキストモダリティニューラルネットワーク20、音声モダリティニューラルネットワーク30及び共通ニューラルネットワーク40を学習する。図3及び4を参照して上述した実施例による学習処理では、隠れ層の重み行列は共有される共通ニューラルネットワーク40において同じとされたが、テキストモダリティニューラルネットワーク20及び音声モダリティニューラルネットワーク30から入力された各ベクトルに対する共通ニューラルネットワーク40の隠れ層から出力されるベクトルが互いに近いものになることを明示的に保証するものでない。このため、共通ニューラルネットワーク40の入力層に隠れ層から出力されるベクトルが近似したものになるよう共通ニューラルネットワーク40を学習することによって、より精度の高い変換が可能になると考えられる。

#### 【0047】

具体的には、学習装置100は、テキストデータと音響特徴量とのペアから構成される訓練データに対して、当該テキストデータをテキストモダリティニューラルネットワーク20に入力し、テキストモダリティニューラルネットワーク20から出力されたベクトルを取得する。そして、学習装置100は、取得したベクトルを共通ニューラルネットワーク40に入力し、共通ニューラルネットワーク40から出力された音響特徴量を取得し、取得した音響特徴量と訓練データの音響特徴量との間の誤差(LOSS<sub>main</sub>)を算出する。

#### 【0048】

一方、学習装置100は、音声波形データと音響特徴量とのペアから構成される訓練データに対して、当該音声波形データを音声モダリティニューラルネットワーク30に入力し、音声モダリティニューラルネットワーク30から出力されたベクトルを取得する。そして、学習装置100は、取得したベクトルを共通ニューラルネットワーク40に入力し

10

20

30

40

50

、共通ニューラルネットワーク40の一部のレイヤ（例えば、入力層からL番目の隠れ層）から構成されるサブニューラルネットワークから出力されたベクトル（ $h^1_{sub}$ ）を取得する一方、テキストモダリティニューラルネットワーク20から共通ニューラルネットワーク40に入力されたベクトルに対して、当該サブニューラルネットワークから出力されたベクトル（ $h^1_{main}$ ）を取得する。

【0049】

その後、学習装置100は、2つのベクトル（ $h^1_{main}$ 、 $h^1_{sub}$ ）の間の距離に基づき誤差（ $loss_{sub}$ ）を算出し、誤差（ $loss_{main}$ ）と誤差（ $loss_{sub}$ ）との加重和に基づきテキストモダリティニューラルネットワーク20、音声モダリティニューラルネットワーク30及び共通ニューラルネットワーク40を学習する。例えば、学習装置100は、

$$loss = loss_{main} + \lambda \cdot distance(h^1_{main}, h^1_{sub})$$

に従って（ $\lambda$ は、スカラー値である）、2つの誤差（ $loss_{main}$ 、 $loss_{sub}$ ）の加重和 $loss$ を算出してもよい。ここで、距離 $distance$ は、例えば、コサイン距離であってもよい。

【0050】

学習装置100は、算出した誤差の加重和が減少するように、例えば、バックプロパゲーションに従ってテキストモダリティニューラルネットワーク20、音声モダリティニューラルネットワーク30及び共通ニューラルネットワーク40のパラメータ（例えば、隠れ層の重み行列）を更新する。

【0051】

図6は、本発明の他の実施例による学習処理を示すフローチャートである。当該学習処理は、学習装置100、具体的には、学習装置100のプロセッサ101によって実行される。

【0052】

図6に示されるように、ステップS201において、学習装置100は、訓練データを取得する。

【0053】

ステップS202において、学習装置100は、処理対象の訓練データがテキストデータと音響特徴量とのペアである場合、ステップS203に進み、処理対象の訓練データが音声波形データと音響特徴量とのペアである場合、ステップS206に進む。

【0054】

ステップS203において、学習装置100は、訓練データのテキストデータをテキストモダリティニューラルネットワーク20に入力し、テキストモダリティニューラルネットワーク20から出力されたベクトルを取得する。

【0055】

ステップS204において、学習装置100は、取得したベクトルを共通ニューラルネットワーク40に入力する。

【0056】

ステップS205において、学習装置100は、共通ニューラルネットワーク40から出力された音響特徴量を取得すると共に、共通ニューラルネットワーク40のサブニューラルネットワーク（例えば、入力層から所定番目の隠れ層）から出力されたベクトル（ $h^1_{main}$ ）を取得する。

【0057】

一方、ステップS206において、学習装置100は、訓練データの音声波形データを音声モダリティニューラルネットワーク30に入力し、音声モダリティニューラルネットワーク30から出力されたベクトルを取得する。

【0058】

ステップS207において、学習装置100は、取得したベクトルを共通ニューラルネ

10

20

30

40

50

ットワーク 40 に入力する。

【0059】

ステップ S 208 において、学習装置 100 は、共通ニューラルネットワーク 40 のサブニューラルネットワークから出力されたベクトル ( $h^1_{sub}$ ) を取得する。

【0060】

ステップ S 209 において、学習装置 100 は、共通ニューラルネットワーク 40 から取得した音響特徴量と訓練データの音響特徴量との誤差 ( $loss_{main}$ ) と、2つのベクトル ( $h^1_{main}$ ,  $h^1_{sub}$ ) の間の距離 ( $loss_{sub}$ ) とを算出する。

【0061】

ステップ S 210 において、学習装置 100 は、ステップ S 209 において算出した誤差 ( $loss_{main}$ ) と距離 ( $loss_{sub}$ ) との加重和 ( $loss$ ) を算出し、算出した加重和 ( $loss$ ) が減少するように、例えば、バックプロパゲーションに従ってテキストモダリティニューラルネットワーク 20、音声モダリティニューラルネットワーク 30 及び共通ニューラルネットワーク 40 のパラメータ (例えば、隠れ層の重み行列) を更新する。

【0062】

学習装置 100 は、所定の終了条件を充足するまで、各訓練データに対して上述したステップ S 201 ~ S 210 を繰り返す。当該所定の終了条件は、例えば、所定の回数の繰り返しを終了したこと、誤差 ( $loss$ ) が所定の閾値以下になったこと、誤差 ( $loss$ ) が収束したことなどであってもよい。

【0063】

なお、2つのタイプのニューラルネットワーク構造の学習処理について個別に説明したが、これら2つのタイプの学習処理が組み合わせ可能であることは当業者に理解されるであろう。この場合、誤差 ( $loss$ ) は、例えば、

$$loss = loss_{main} + loss_{sub} + l^L distance(h^1_{main}, h^1_{sub})$$

に従って算出され、テキストモダリティニューラルネットワーク 20、音声モダリティニューラルネットワーク 30 及び共通ニューラルネットワーク 40 のパラメータが、誤差を減少させるように更新されると共に、2つの共通ニューラルネットワーク 40 のパラメータが同期的に学習される。

[ 共通ニューラルネットワーク 40 に対する話者適応処理 ]

次に、図 7 及び 8 を参照して、本発明の一実施例による共通ニューラルネットワーク 40 に対する話者適応処理を説明する。本実施例では、上述した学習処理に従ってニューラルネットワーク構造 10 を学習した後、所与の未知話者の訓練データが与えられると、学習装置 100 は、当該訓練データに応じて、テキストモダリティニューラルネットワーク 20 及び共通ニューラルネットワーク 40 と、音声モダリティニューラルネットワーク及び共通ニューラルネットワーク 40 とを選択的に利用して、共通ニューラルネットワーク 40 の話者空間における当該未知話者を示す話者コードベクトルを推定する。

【0064】

図 7 は、本発明の一実施例による未知話者適応処理を示す概略図である。本実施例では、図 7 に示されるように、与えられた訓練データが所与の未知話者のテキストデータと音響特徴量とのペアである場合、学習装置 100 は、テキストモダリティニューラルネットワーク 20 及び共通ニューラルネットワーク 40 を利用して、当該未知話者の話者コードベクトルを推定する。他方、与えられた訓練データが所与の未知話者の音声波形データと音響特徴量とのペアである場合、学習装置 100 は、音声モダリティニューラルネットワーク 30 及び共通ニューラルネットワーク 40 を利用して、当該未知話者の話者コードベクトルを推定する。

【0065】

具体的には、学習装置 100 は、所与の未知話者の訓練データがテキストデータと音響特徴量とから構成される場合、当該テキストデータをテキストモダリティニューラルネッ

10

20

30

40

50

トワーク 20 に入力し、テキストモダリティニューラルネットワーク 20 から取得したベクトルを共通ニューラルネットワーク 40 に入力し、共通ニューラルネットワーク 40 から取得した音響特徴量と訓練データの音響特徴量との間の誤差に基づき当該話者の話者コードベクトルを決定する。他方、学習装置 100 は、所与の未知話者の訓練データが音声波形データと音響特徴量とから構成される場合、音声波形データを音声モダリティニューラルネットワーク 30 に入力し、音声モダリティニューラルネットワーク 30 から取得したベクトルを共通ニューラルネットワーク 40 に入力し、共通ニューラルネットワーク 40 から取得した音響特徴量と訓練データの音響特徴量との間の誤差に基づき当該話者の話者コードベクトルを決定する。

【0066】

例えば、図 1 に示される具体例によると、話者コードベクトル  $d^{(i)}$  は、

$$d^{(i)} = d^{(i)} + W_D^T f_{n-1}$$

に従って更新される。ここで、 $f_{n-1}$  は所定値以下の小さな値であり、 $f$  は誤差伝搬のための関数であり、

$$f_{n-1}(c) = W_{C,N}(c) \cdot \tau^{-1}(e')$$

として定義され、 $\tau^{-1}$  は活性化関数によって決定される伝搬用の関数であり、 $e'$  は共通ニューラルネットワーク 40 から取得した音響特徴量と訓練データの音響特徴量との間の誤差の微分値である。なお、当該未知話者適応処理では、共通ニューラルネットワーク 40 の重み行列  $W$  及びバイアスベクトル  $b$  は更新されない。

【0067】

このようにして、共通ニューラルネットワーク 40 における話者コードベクトル（潜在変数）を特定することによって、学習済みのニューラルネットワーク構造 10 を特定の未知話者に適応させることができる。

【0068】

図 8 は、本発明の一実施例による未知話者適応処理を示すフローチャートである。当該学習処理は、学習装置 100、具体的には、学習装置 100 のプロセッサ 101 によって実行される。

【0069】

図 8 に示されるように、ステップ S301 において、学習装置 100 は、所与の未知話者の訓練データを取得する。

【0070】

ステップ S302 において、学習装置 100 は、訓練データがテキストデータと音響特徴量とのペア又は音声波形データと音響特徴量とのペアから構成されているか判断し、訓練データがテキストデータと音響特徴量とのペアから構成されている場合、ステップ S303 に進み、訓練データが音声波形データと音響特徴量とのペアから構成されている場合、ステップ S306 に進む。

【0071】

ステップ S303 において、学習装置 100 は、訓練データのテキストデータをテキストモダリティニューラルネットワーク 20 に入力し、テキストモダリティニューラルネットワーク 20 から出力されたベクトルを取得する。

【0072】

ステップ S304 において、学習装置 100 は、取得したベクトルを共通ニューラルネットワーク 40 に入力し、共通ニューラルネットワーク 40 から出力された音響特徴量を取得する。

【0073】

ステップ S305 において、学習装置 100 は、共通ニューラルネットワーク 40 から取得した音響特徴量と訓練データの音響特徴量との間の誤差を算出する。

【0074】

一方、ステップ S306 において、学習装置 100 は、訓練データの音声波形データを音声モダリティニューラルネットワーク 30 に入力し、音声モダリティニューラルネット

10

20

30

40

50

ワーク 30 から出力されたベクトルを取得する。

【0075】

ステップ S 307 において、学習装置 100 は、取得したベクトルを共通ニューラルネットワーク 40 に入力し、共通ニューラルネットワーク 40 から出力された音響特徴量を取得する。

【0076】

ステップ S 308 において、学習装置 100 は、共通ニューラルネットワーク 40 から取得した音響特徴量と訓練データの音響特徴量との間の誤差を算出する。

【0077】

ステップ S 309 において、学習装置 100 は、ステップ S 305 及び S 308 において算出した誤差が減少するように、例えば、上述した更新式を利用してバックプロパゲーションに従って共通ニューラルネットワーク 40 の話者コードベクトルを更新する。

【0078】

学習装置 100 は、所定の終了条件を充足するまで、各訓練データに対して上述したステップ S 301 ~ S 309 を繰り返す。当該所定の終了条件は、例えば、所定の回数の繰り返しを終了したこと、誤差が所定の閾値以下になったこと、誤差が収束したことなどであってもよい。

[ 学習済みニューラルネットワーク構造を利用した音声合成処理 ]

次に、図 9 ~ 11 を参照して、本発明の一実施例による音声合成処理を説明する。本実施例では、音声合成装置 200 は、上述した学習装置 100 によって特定の話者に対して学習されたテキストモダリティニューラルネットワーク 20 及び共通ニューラルネットワーク 40 を利用して、音声合成対象のテキストデータから当該話者に対応する音声データを生成及び再生する。

【0079】

図 9 は、本発明の一実施例による音声合成処理を示す概略図である。本実施例では、音声合成装置 200 は、音声合成対象のテキストデータが与えられると、図 9 に示されるように、上述した学習装置 100 によって特定の話者に対して学習されたテキストモダリティニューラルネットワーク 20 及び共通ニューラルネットワーク 40 を利用して、当該テキストデータから当該話者に対応する音響特徴量を生成する。具体的には、音声合成装置 200 は、入出力インタフェース 104 を介して、テキストデータを取得し、当該話者に対応するテキストデータから生成された音響特徴量を再生してもよい。

【0080】

図 10 は、本発明の一実施例による音声合成処理を示すフローチャートである。当該音声合成処理は、音声合成装置 200、具体的には、音声合成装置 200 のプロセッサ 101 によって実行される。

【0081】

図 10 に示されるように、ステップ S 401 において、音声合成装置 200 は、音声合成対象となるテキストデータを取得する。例えば、テキストデータは、音声合成装置 200 の入出力インタフェース 104 を介し入力されたものであってもよい。

【0082】

ステップ S 402 において、音声合成装置 200 は、取得したテキストデータを学習済みテキストモダリティニューラルネットワーク 20 に入力し、テキストモダリティニューラルネットワーク 20 から出力されたベクトルを取得する。

【0083】

ステップ S 403 において、音声合成装置 200 は、取得したベクトルを学習済み共通ニューラルネットワーク 40 に入力し、共通ニューラルネットワーク 40 から出力された音響特徴量を取得する。

【0084】

ステップ S 404 において、音声合成装置 200 は、共通ニューラルネットワーク 40 から取得した特定の話者に対応する音響特徴量を何れかの音声データフォーマットに変換

10

20

30

40

50

し、変換された音声データを再生する。例えば、変換された音声データは、当該話者の声、テンポ、アクセントなどに近い音声によって入力されたテキストデータを再生したものとなりうる。

【0085】

図11は、本発明の各種実施例による学習処理の実験結果を示す図である。図11において、VL、SS、JG、TL及びJG+TLは、上述した学習済みニューラルネットワーク構造を利用したものを含む各種音声合成システムを表す。

【0086】

VLは、3つのニューラルネットワークから構成されるニューラルネットワーク構造10でなく、従来のニューラルネットワーク構造を利用したシステムである。SSは、ニューラルネットワーク構造10の各モダリティニューラルネットワークを単純に置き換えて学習されたシステムである。JGは、図3及び4を参照して説明した学習処理により学習されたニューラルネットワーク構造10を利用したシステムである。TLは、図5及び6を参照して説明した学習処理により学習されたニューラルネットワーク構造10を利用したシステムである。JG+TLは、JGとTLとを組み合わせた学習処理により学習されたニューラルネットワーク構造10を利用したシステムである。

【0087】

図11では、10、40、160及び320個の未知話者の訓練データによって学習された場合の各音声合成システムの誤差(MCD)のシミュレーション結果が示される。図から理解されうるように、訓練データとして音声データとテキストデータとが与えられる教師有り学習と、音声データのみが与えられる教師なし学習との何れのケースでも、上述した実施例によるJG、TL及びJG+TLは、VL及びSSに対して有意に誤差を低減するという結果を得ることができた。

【0088】

なお、上述した実施例では、テキストデータと音声データとが異なるモダリティとして扱われたが、本発明は、これに限定されるものでなく、他のタイプのモダリティの組み合わせに同様に適用可能であることは理解されるであろう。

【0089】

以上、本発明の実施例について詳述したが、本発明は上述した特定の実施形態に限定されるものではなく、特許請求の範囲に記載された本発明の要旨の範囲内において、種々の変形・変更が可能である。

【符号の説明】

【0090】

- 10 ニューラルネットワーク構造
- 20 テキストモダリティニューラルネットワーク
- 30 音声モダリティニューラルネットワーク
- 40 共通ニューラルネットワーク
- 100 学習装置
- 200 音声合成装置

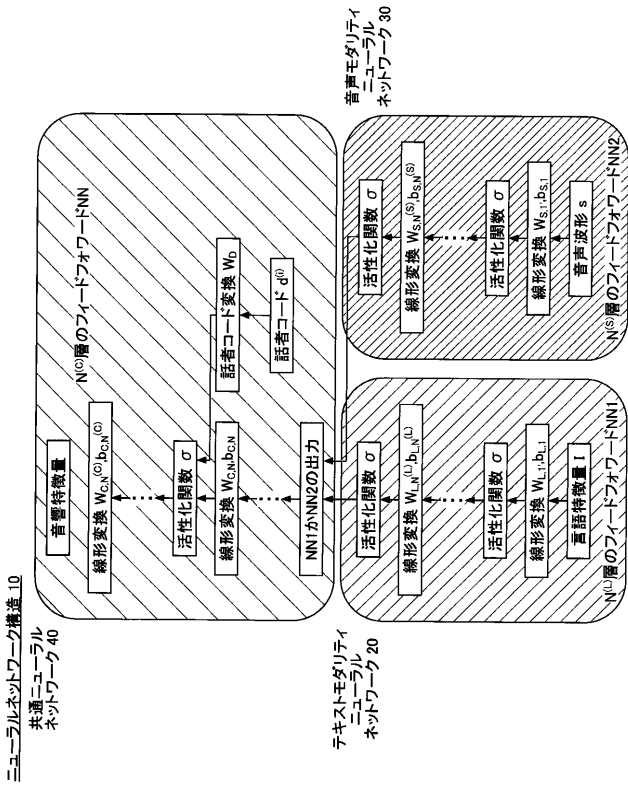
10

20

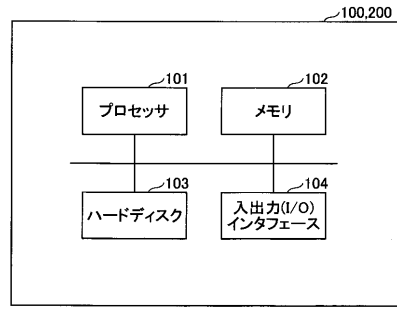
30



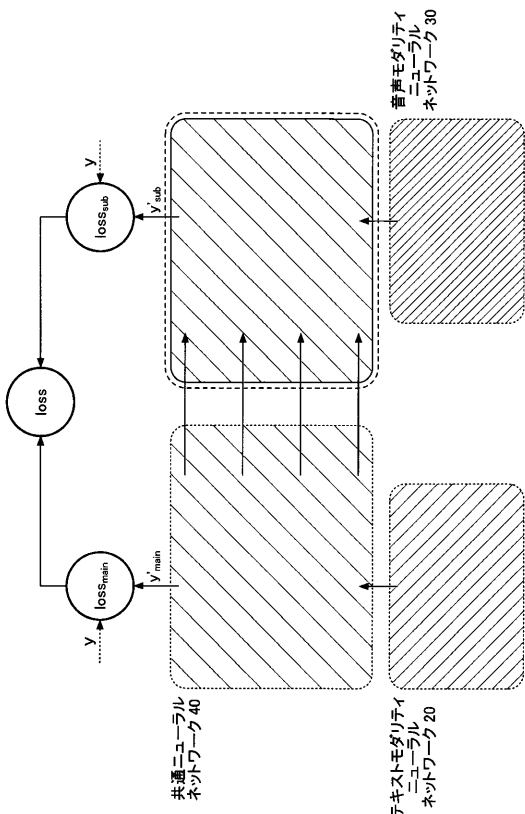
【図1】



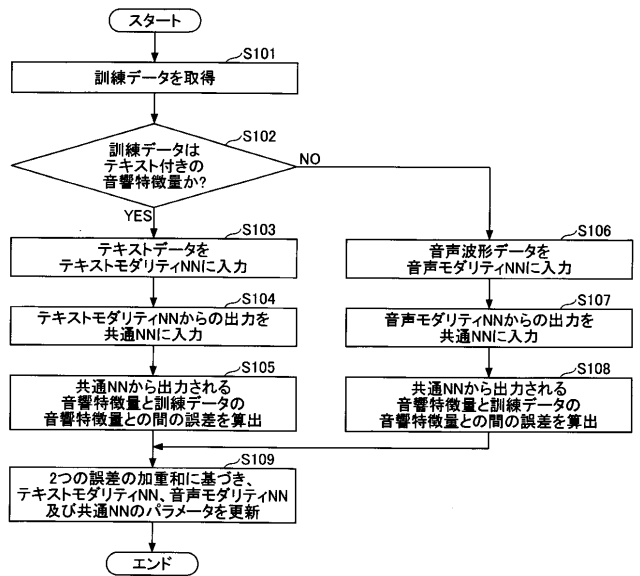
【図2】



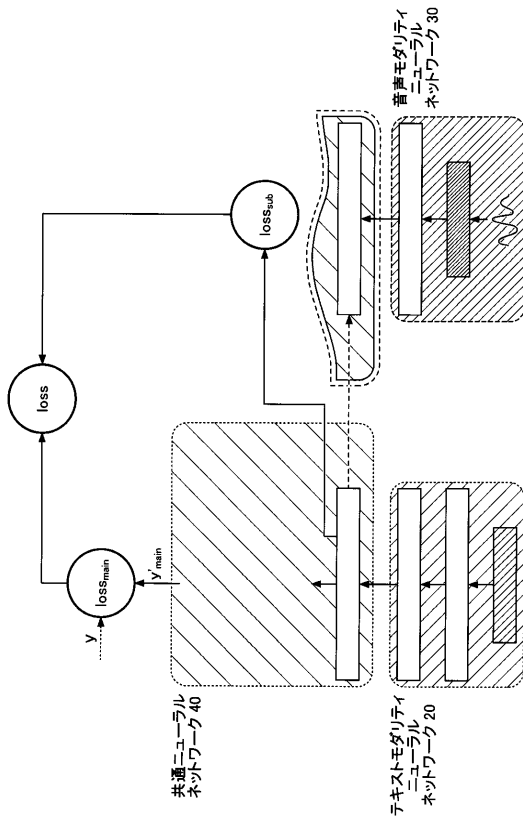
【図3】



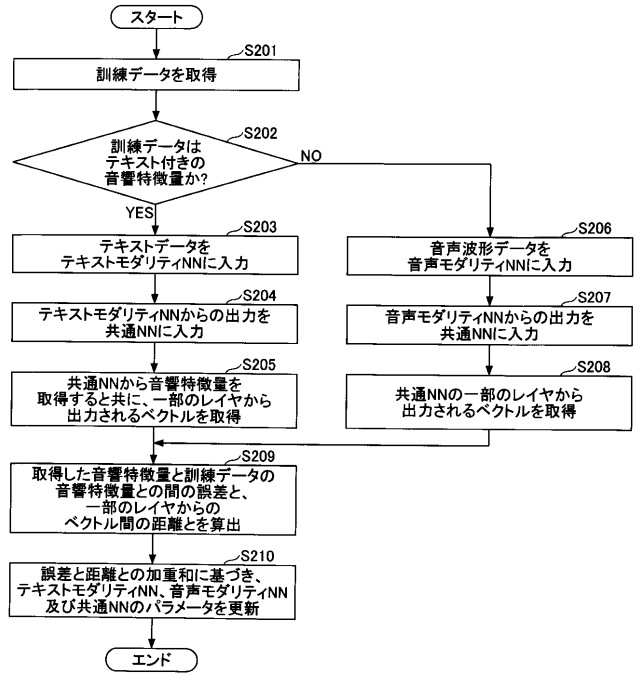
【図4】



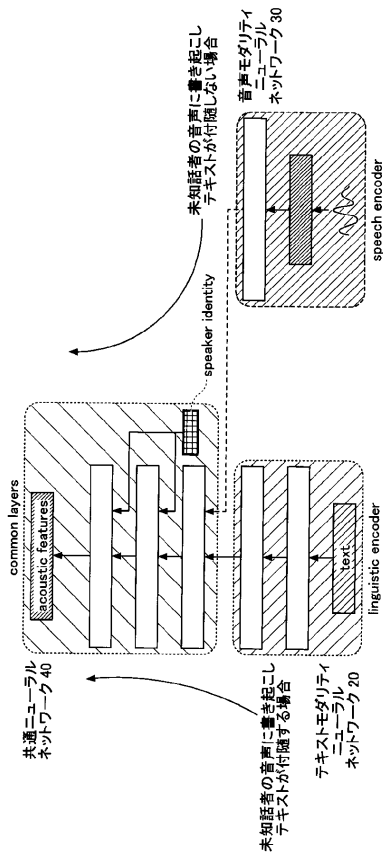
【図5】



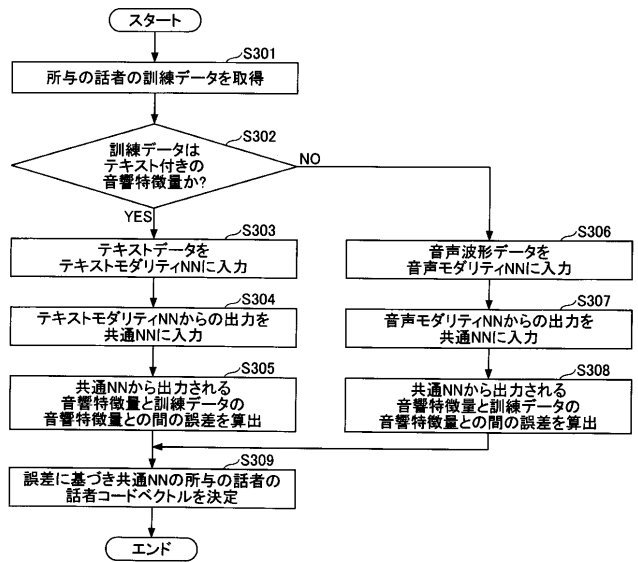
【図6】



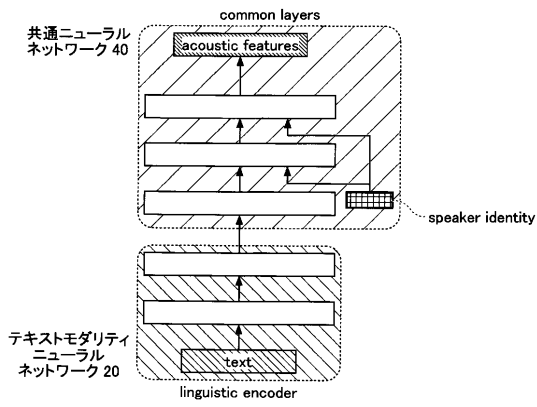
【図7】



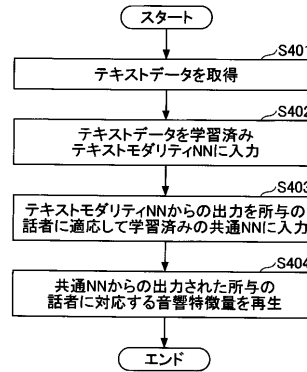
【図8】



【 図 9 】



【 図 1 0 】



【 図 1 1 】

